

CF0 12088 US
08/863,047

日 本 国 特 許 庁
PATENT OFFICE
JAPANESE GOVERNMENT

別紙添付の書類に記載されている事項は下記の出願書類に記載されて
いる事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed
with this Office.

出 願 年 月 日
Date of Application:

1996年 9月 5日

出 願 番 号
Application Number:

平成 8年特許願第232969号

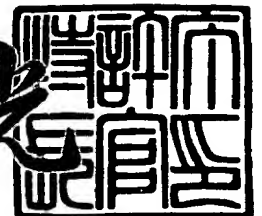
出 願 人
Applicant(s):

キヤノン株式会社

1997年 6月20日

特許庁長官
Commissioner,
Patent Office

荒井 寿光



出証番号 出証特平09-3046808

【書類名】 特許願

【整理番号】 3328019

【提出日】 平成 8年 9月 3日

【あて先】 特許庁長官 荒井 寿光 殿

【国際特許分類】 G06F 15/20

【発明の名称】 情報検索方法及びその装置、記憶媒体

【請求項の数】 7

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社
内

【氏名】 伊藤 史朗

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社
内

【氏名】 大谷 紀子

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社
内

【氏名】 柴田 昇吾

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社
内

【氏名】 上田 隆也

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社
内

【氏名】 池田 裕治

【特許出願人】

【識別番号】 000001007

【郵便番号】 146
【住所又は居所】 東京都大田区下丸子3丁目30番2号
【氏名又は名称】 キヤノン株式会社
【代表者】 御手洗 富士夫
【電話番号】 03-3758-2111

【代理人】

【識別番号】 100069877
【郵便番号】 146
【住所又は居所】 東京都大田区下丸子3丁目30番2号キヤノン株式会社
内

【弁理士】
【氏名又は名称】 丸島 儀一
【電話番号】 03-3758-2111

【手数料の表示】

【予納台帳番号】 011224
【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1
【物件名】 図面 1
【物件名】 要約書 1
【包括委任状番号】 9003707

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 情報検索方法及びその装置、記憶媒体

【特許請求の範囲】

【請求項1】 情報を集合として保持する情報保持手段と、
検索条件を保持する検索条件保持手段と、
前記情報保持手段に保持されている各要素から、前記検索条件保持手段に保持されている検索条件を満足する要素を検索する検索手段と、
前記検索条件の検索結果を用いて、前記情報保持手段に保持されている各情報集合のスコアを演算する演算手段と、
前記演算手段により演算されたスコアを保持する集合スコア保持手段と、
を有することを特徴とする情報検索装置。

【請求項2】 前記演算手段は、集合の要素数と集合中の検索条件を満足する要素数とから集合全体のスコアの統計的区間推定を行いその下限値をもってスコアとすることを特徴とする請求項1に記載の情報検索装置。

【請求項3】 前記演算手段は、集合の要素数と各要素の検索条件に対するスコアを用いて集合全体のスコアの統計的区間推定を行いその下限値をもってスコアとすることを特徴とする請求項1に記載の情報検索装置。

【請求項4】 情報を集合として保持する情報保持工程と、
検索条件を保持する検索条件保持工程と、
前記情報保持工程により保持されている各要素から、前記検索条件保持工程により保持されている検索条件を満足する要素を検索する検索工程と、
前記検索工程の検索結果を用いて、前記情報保持工程により保持されている各情報集合のスコアを演算する演算工程と、
前記演算工程により演算されたスコアを保持する集合スコア保持工程と、
を有することを特徴とする情報検索方法。

【請求項5】 前記演算工程は、集合の要素数と集合中の検索条件を満足する要素数とから集合全体のスコアの統計的区間推定を行いその下限値をもってスコアとすることを特徴とする請求項4に記載の情報検索方法。

【請求項6】 前記演算工程は、集合の要素数と各要素の検索条件に対するスコアを用いて集合全体のスコアの統計的区間推定を行いその下限値をもってスコアとすることを特徴とする請求項1に記載の情報検索方法。

【請求項7】 情報を集合として保持する情報保持工程と、
検索条件を保持する検索条件保持工程と、
前記情報保持工程により保持されている各要素から、前記検索条件保持工程により保持されている検索条件を満足する要素を検索する検索工程と、

前記検索工程の検索結果を用いて、前記情報保持工程により保持されている各情報集合のスコアを演算する演算工程と、

前記演算工程により演算されたスコアを保持する集合スコア保持工程と、
を有することを特徴とする情報検索プログラムを格納したコンピュータで読み取り可能な記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、情報の集合を検索する情報検索方法及びその装置に関するものである。

【0002】

【従来の技術】

情報を含む文書やレコードなどを集合として保持し、情報を利用する場合には集合として検索する装置がある。例えば、従来の文書処理装置（特願平8-129899）では、文書を集合として保持するものとしてフォルダを用意し、フォルダ単位で検索を行なっている。

【0003】

従来、この種の文書処理装置では、検索条件に対する集合のスコアとして検索条件を満足する文書やレコードが集合中に含まれる割合を採用していた。この割合が高いほど、集合として検索条件に合う可能性が高いと考えられる。

【0004】

【発明が解決しようとしている課題】

しかしながら、上記従来例の装置では、集合中の要素数の多寡がスコアに反映されないため、要素数が異なる集合同士の比較がしづらいという問題点がある。例えば、集合の要素数が1で条件を満足する要素が1の場合と、集合の要素数が10で条件を満足する要素が9の場合では、前者のスコアが高くなるが、集合として検索条件に合う可能性が高いのは後者である。

【0005】

本発明は上記の問題に鑑みてなされたものであり、集合の要素数の多寡を反映して、集合のスコアを算出する情報検索方法及び装置を提供することを目的とする。

【0006】

【課題を解決する為の手段】

上記課題を解決する為に、本発明は、情報を集合として保持する情報保持手段と、検索条件を保持する検索条件保持手段と、前記情報保持手段に保持されている各要素から、前記検索条件保持手段に保持されている検索条件を満足する要素を検索する検索手段と、前記検索手段の検索結果を用いて、前記情報保持手段に保持されている各情報集合のスコアを演算する演算手段と、前記演算手段により演算されたスコアを保持する集合スコア保持手段とを有することを特徴とする情報検索装置を提供する。

【0007】

【発明の実施の形態】

以下、図面を参照して本発明の実施形態を詳細に説明する。

【0008】

図1は、本発明の一実施形態に係る装置の基本構成を示すブロック図である。

【0009】

同図において、101は検索対象文書を保持する文書保持部である。102は文書の集合を保持する文書集合保持部である。103は検索条件を保持する検索条件保持部である。104は文書保持部101に保持されている文書から、検索

条件保持部103に保持されている検索条件を満足する文書を検索する文書検索部である。105は文書検索部104の検索結果を保持する検索結果保持部である。106は検索結果保持部105に保持されている検索結果を用いて、文書集合保持部102に保持されている各文書集合のスコアを計算する集合スコア演算部である。107は集合スコア演算部106により計算されたスコアを保持する集合スコア保持部である。

【0010】

図2は、本実施形態の情報検索装置の具体的構成を示す図である。

【0011】

同図において、201はCPUであり、後述する手順を実現するプログラムに従って動作する。202はRAMであり、検索条件保持部103と検索結果保持部105と集合スコア保持部107と上記プログラムの動作に必要な記憶領域とを提供する。203はROMであり、後述する手順を実現するプログラムを保持する。204はディスク装置であり、文書保持部101と文書集合保持部102を実現する。205はバスである。206は後述するプログラムにより検索された検索結果を表示するCRTである。207は検索条件等を入力するキーボード(KB)である。

【0012】

なお、本実施例では、文書集合保持部102は、文書集合に一意の集合番号を付与して、集合番号が示す文書集合に属する文書の文書番号のリストを保持する構造になっている。図4に文書集合保持部の例を示す。列401には集合番号が保持され、列402には文書番号のリストが保持されている。

【0013】

この他、文書保持部101は、文書に一意の文書番号を付与して、文書番号が示す文書本文を保持する構造になっている。検索条件保持部103は、検索語のリストを保持する構造になっている。検索結果保持部105は、文書番号のリストを保持する構造になっている。集合スコア保持部107は、集合番号が示す文書集合のスコアを保持する構造になっている。

【0014】

以下、図3のフローチャートを参照して、本実施形態の情報検索装置における処理の手順を示す。

【0015】

ステップS301では、検索条件保持部103に検索語のリストから成る検索条件cが保持されているか否かを調べ、保持されている場合はステップS302に移る。保持されていない場合は、ステップS301を繰り返す。

【0016】

ステップS302では、文書保持部101中の文書から、検索条件保持部103に保持されている検索条件cを満足する文書を検索する。文書の本文中にc中の各語が出現するか否かをパターンマッチングにより調べて、全ての語が出現する場合に、その文書は検索条件cを満足すると判断する。検索条件を満足する文書の文書番号を検索結果保持部105に保持する。そして、ステップS303に移る。

【0017】

ステップS303では、kの値を1に設定する。そして、ステップS304に移る。

【0018】

ステップS304では、kの値と文書集合保持部102中の集合数Nを比較し、 $k \leq N$ であればステップS305に移る。 $k > N$ であれば処理を終了する。

【0019】

ステップS305では、文書集合保持部102中のk番目の文書集合 D_k のスコア s_k を計算する。まず、文書集合 D_k に属する文書数をnとし、 D_k に属する文書のうち、検索結果保持部105中にも含まれる文書の数をもとにする。その上で、自由度 (ϕ_1, ϕ_2) のF分布を用いて、次のように s_k を計算する。

【0020】

【外1】

$$s_k = \frac{\phi_2}{\phi_1 F_{\phi_1}^{\phi_2}(\alpha) + \phi_2}$$

ここで、自由度は、 $\phi_1 = 2(n - x + 1)$ 、 $\phi_2 = 2x$ である。

【0021】

また、 α は区間推定における信頼度を指定するパラメータであるが、例えば0.1とする。そして、ステップS306に移る。

【0022】

ステップS306では、ステップS305で求めたスコア s_k を集合スコア保持部107に保持する。そして、ステップS304に戻る。

【0023】

例えば、図4の文書集合保持部の例において、ステップS302が終了した後、検索結果が(1, 3, 5)となったとすると、文書集合ごとの n と x の値は、 D_1 では $n=5$ 、 $x=3$ 、 D_2 では $n=1$ 、 $x=1$ 、 D_3 では $n=3$ 、 $x=1$ となる。従って、各文書集合のスコア s_k は、

【0024】

【外2】

$$s_1 = \frac{6}{6F_6^6(0.1)+6} \approx 0.25$$

$$s_2 = \frac{2}{2F_2^2(0.1)+2} \approx 0.10$$

$$s_3 = \frac{2}{6F_6^2(0.1)+2} \approx 0.03$$

となる。

【0025】

このように本実施形態の情報検索装置を用いると、検索条件により合う文書集合に対して高いスコアが与えられるので、得られたスコアを用いると、利用者が検索条件に合う文書集合を探すことが容易になる。

【0026】

(第2の実施形態)

上記実施形態においては、集合の要素数と集合中で検索条件を満足する要素数とから2項分布の統計的区間推定を行ないその下限値をもって集合全体のスコア

とする場合について説明した。

【0027】

本実施形態では、集合の要素数と各要素の検索条件に対するスコアを用いて母平均の区間推定を行ないその下限値をもって集合全体のスコアとする場合について説明する。

【0028】

なお、本実施形態の基本構成は図1と同様である。ただし、文書検索部104では、各文書の検索条件に対するスコアを算出し、検索結果保持部105では、文書ごとのスコアを保持する。図6に検索結果保持部105の例を示す。列601は文書番号を保持し、列602がその文書のスコアを保持する。

【0029】

以下、図5のフローチャートを参照して、本実施形態の情報検索装置における処理の手順を示す。

【0030】

ステップS501では、検索条件保持部103に検索語のリストから成る検索条件cが保持されているか否かを調べ、保持されている場合はステップS502に移る。保持されていない場合は、ステップS501を繰り返す。

【0031】

ステップS502では、文書保持部101中の文書について、検索条件保持部103に保持されている検索条件cに対するスコアを計算する。スコアは、文書の本文中にc中の各語が出現する頻度を用いて一般的に計算される。計算されたスコアを検索結果保持部105に保持する。そして、ステップS503に移る。

【0032】

ステップS503では、kの値を1に設定する。そして、ステップS504に移る。

【0033】

ステップS504では、kの値と文書集合保持部102中の集合数Nを比較し、 $k \leq N$ であればステップS505に移る。 $k > N$ であれば処理を終了する。

【0034】

ステップS505では、文書集合保持部102中のk番目の集合 D_k に属する文書数を n とする。また、 D_k に属する文書のスコアの平均を x とする。 $n > 1$ の場合は、不偏推定量 V を次の式により求め、

【0035】

【外3】

$$V = \frac{\sum (x - \bar{x})^2}{n-1}$$

その上で、自由度 ϕ 、両側確率 α の t 分布を用いて、次のように s_k を計算する。

【0036】

【外4】

$$s_k = \bar{x} - t(\phi, \alpha) \frac{\sqrt{V}}{\sqrt{n}}$$

ここで、自由度は、 $\phi = n - 1$ である。 $n = 1$ の場合は、

【0037】

【外5】

$$s_k = \alpha \bar{x}$$

とする。なお、 α は区間推定における信頼度を指定するパラメータであるが、例えば0.1とする。そして、ステップS506に移る。

【0038】

ステップS506では、ステップS505で求めたスコア s を集合スコア保持部107に保持する。そして、ステップS504に戻る。

【0039】

上記実施形態においては、検索条件として検索語のAND条件をとる場合について説明したが、これに限定されるものではない。他の論理関係や文書中の位置条件など、文書に対する任意の検索条件を用いてよい。

【0040】

また、パターンマッチングにより文書の検索を行なう場合について説明したが、これに限定されるものではない。文書に対するインデックスを作成しておき、インデックスを用いて検索するなど、任意の検索手段を用いてよい。

【0041】

また、集合を構成する情報が文書である場合について説明したが、これに限定されるものではない。データの集合であるレコードなど任意のものでよい。この場合、情報に応じた検索手段を用いる。

【0042】

また、各集合についてスコアを計算する場合について説明したが、これに限定されるものではない。文書ごとに、その文書が属する集合を保持しておき、検索結果に含まれる文書を最低1個含む集合についてだけスコアを計算してもよい。それ以外の集合のスコアは0になる。

【0043】

また、集合に対して全てのスコアを保持する場合について説明したが、これに限定されるものではなく、一部のスコアだけを保持してもよい。例えば、スコアが予め設定された閾値を越えたものだけを保持する方法や予め設定された数や割合の範囲内で結果を保持する方法がある。

【0044】

また、各部を同一の計算機上で構成する場合について説明したが、これに限定されるものではなく、ネットワーク上に分散した計算機や処理装置などに分かれて各部を構成してもよい。

【0045】

また、図2は、本実施形態の情報検索装置の具体的構成を示す図である。

【0046】

検索条件保持部と検索結果保持部と集合スコア保持部をRAMで、文書保持部と文書集合保持部をディスク装置で実現する場合について説明したが、これに限定されるものではなく、任意の記憶媒体を用いて実現してもよい。

【0047】

また、プログラムをROMに保持する場合について説明したが、これに限定されるものではなく、任意の記憶媒体を用いて実現してもよい。また、同様の動作をする回路で実現してもよい。

【0048】

また、本発明は、前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ（またはCPUやMPU）が記憶媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。

【0049】

この場合、記憶媒体から読み出されたプログラムコード自体が本発明の新規な機能を実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。

【0050】

プログラムコードを供給するための記憶媒体としては、例えば、フロッピーディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ不揮発性のメモ리카ード、ROMなどを用いることができる。

【0051】

また、コンピュータが読み出したプログラムコードを実行することによって、前述した実施形態の機能が実現される他、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOSなどが実際の処理の一部または全部を行い、その処理によっても前述した実施形態の機能が実現され得る。

【0052】

さらに、記憶媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行い、その処理によっても前述した実施形態の機能が実現され得る。

【0053】

尚、本発明は、前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体からそのプログラムをパソコン通信等通信ラインを介して要求者にそのプログラムを配信する場合にも適用できることは言うまでもない。

【0054】

【発明の効果】

以上説明したように、本発明によれば、集合として保持されている情報に対して集合の要素数の多寡を反映して集合のスコアを算出するようにしたので、検索条件により合う文書集合に対して高いスコアが与えられ、利用者が検索条件に合う文書集合を探すことが容易になるという効果が得られる。

【図面の簡単な説明】

【図1】

本発明に係る情報検索装置の実施形態の基本構成を示すブロック図である。

【図2】

本発明の実施形態の具体的構成を示す図である。

【図3】

本発明の実施形態における処理の概要を示すフローチャートである。

【図4】

本発明の実施形態における文書集合保持部の例を示す図である。

【図5】

本発明の第2の実施形態における処理の概要を示すフローチャートである。

【図6】

本発明の第2の実施形態における検索結果保持部の例を示す図である。

【図7】

本発明の制御プログラムをコンピュータにロードする様子を示す図である。

【符号の説明】

201 CPU

202 RAM

203 ROM

204 ディスク装置

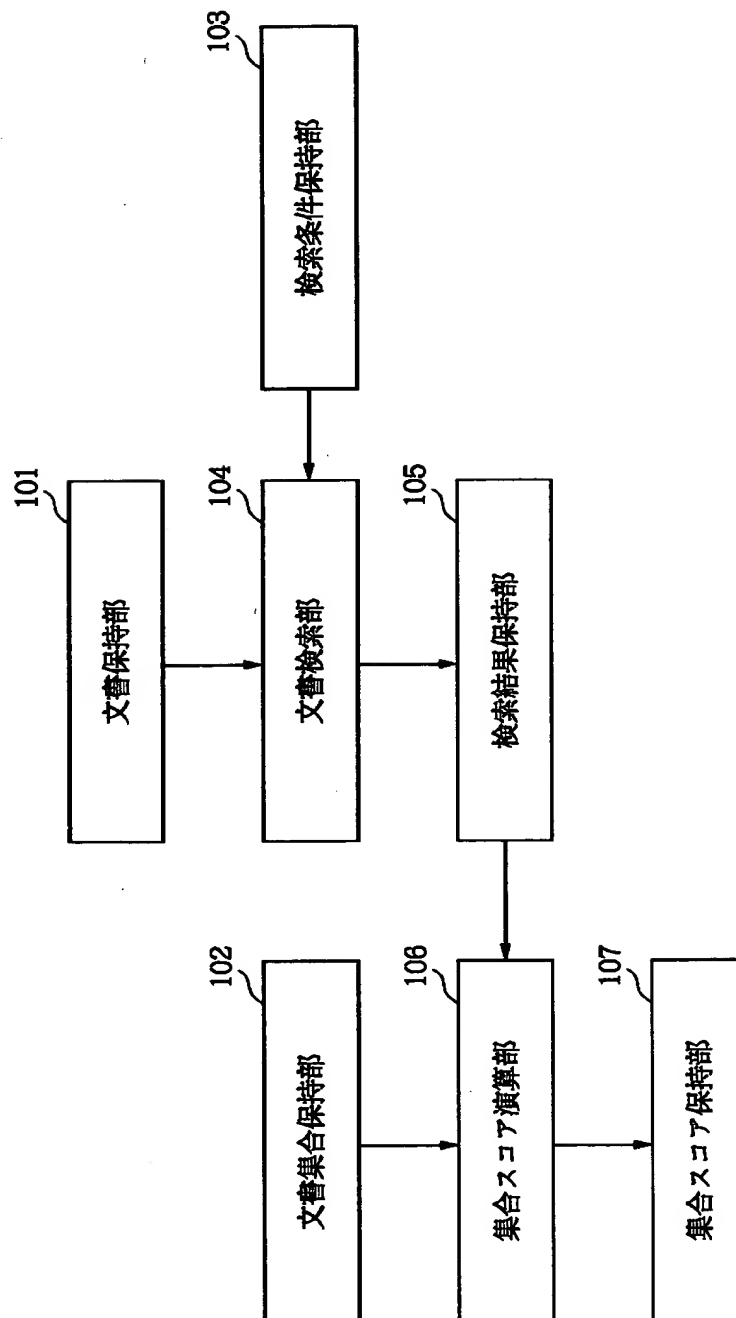
205 バス

206 CRT

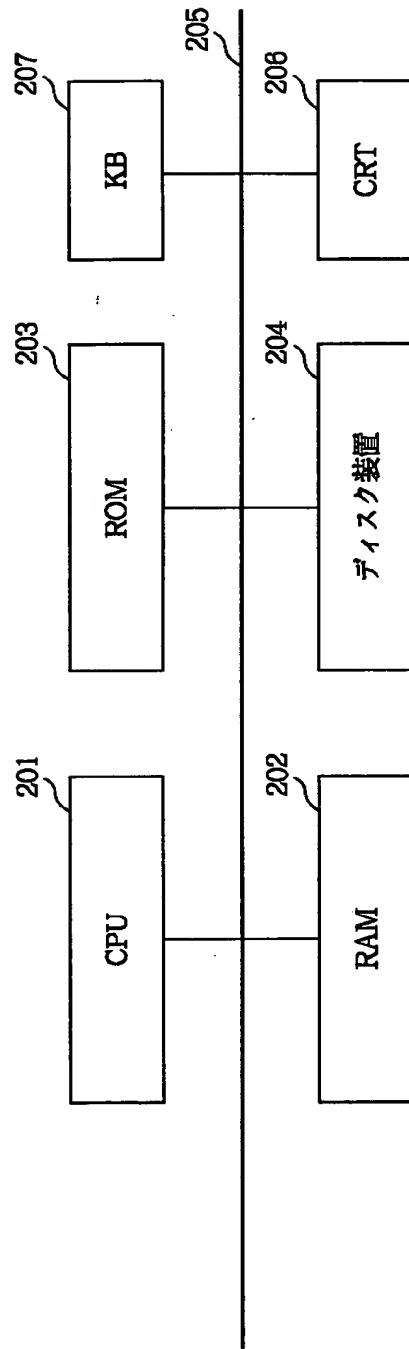
207 KB

【書類名】 図面

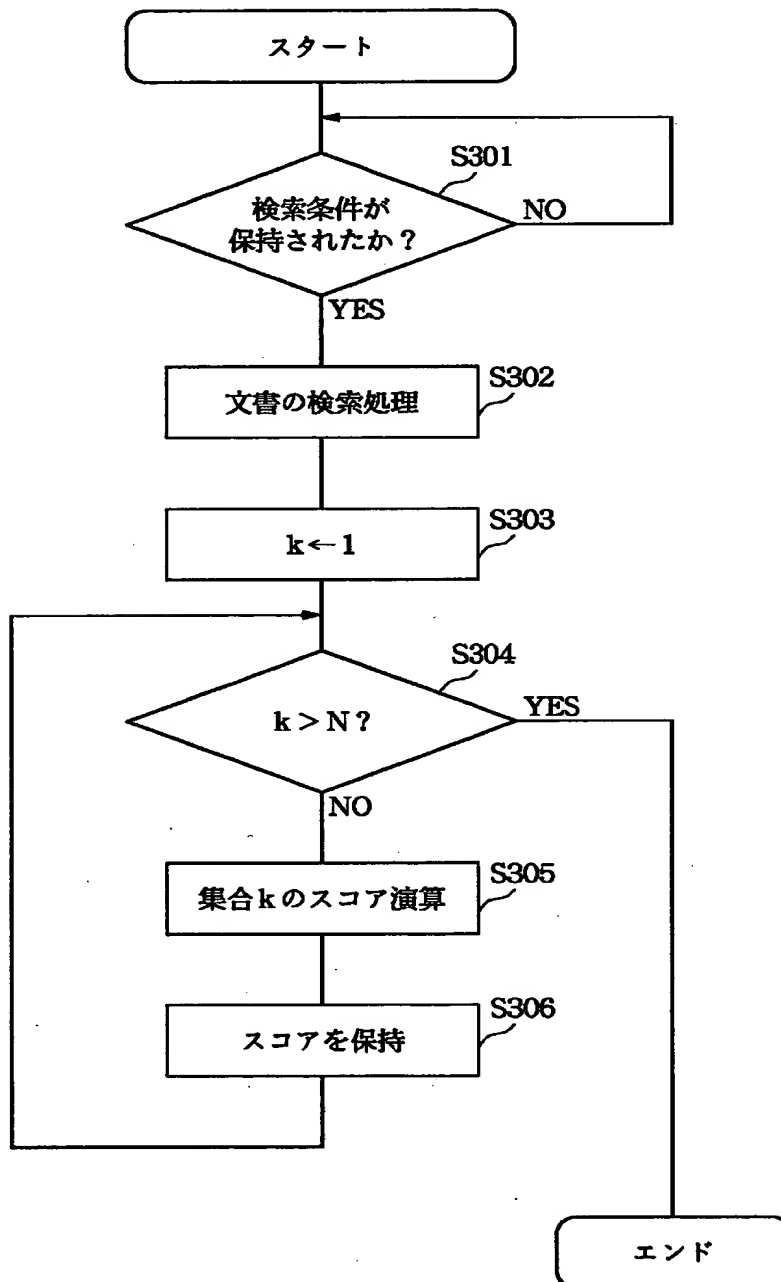
【図1】



【図2】



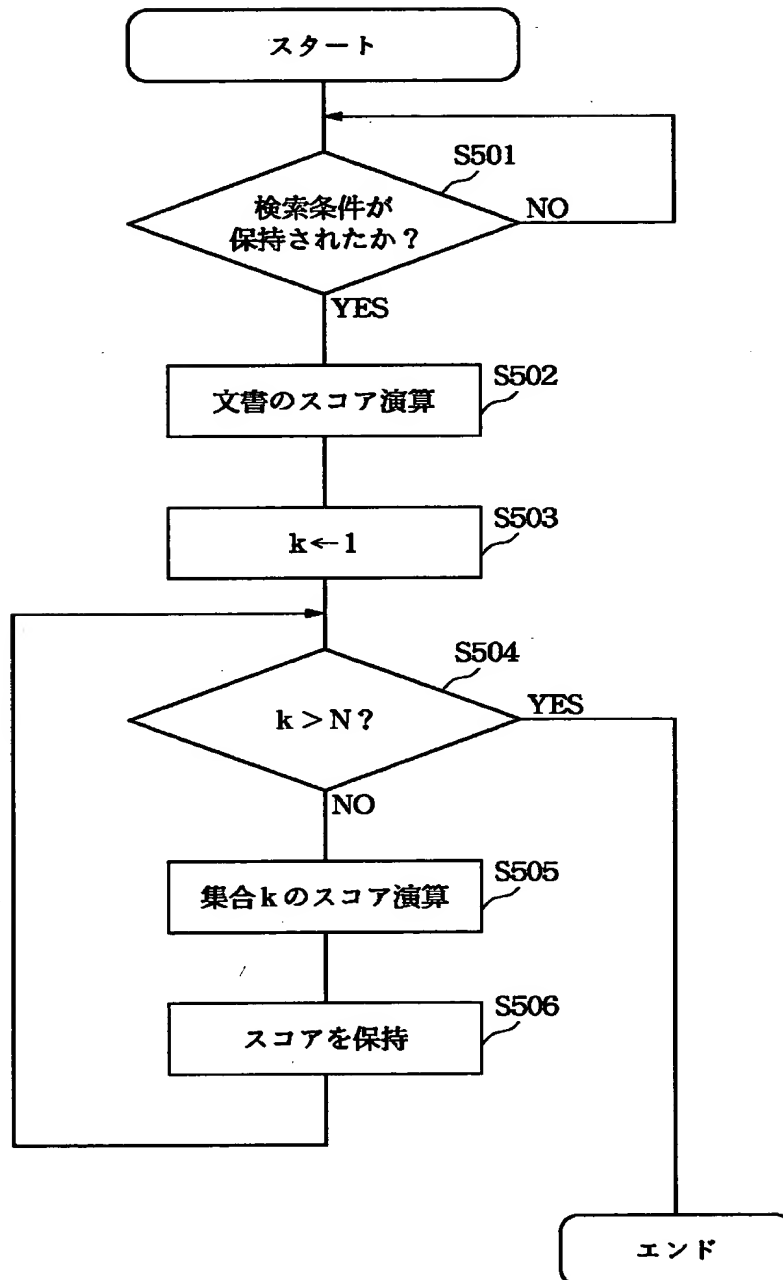
【図3】



【図4】

1	1, 2, 3, 4, 5
2	3
3	2, 5, 6

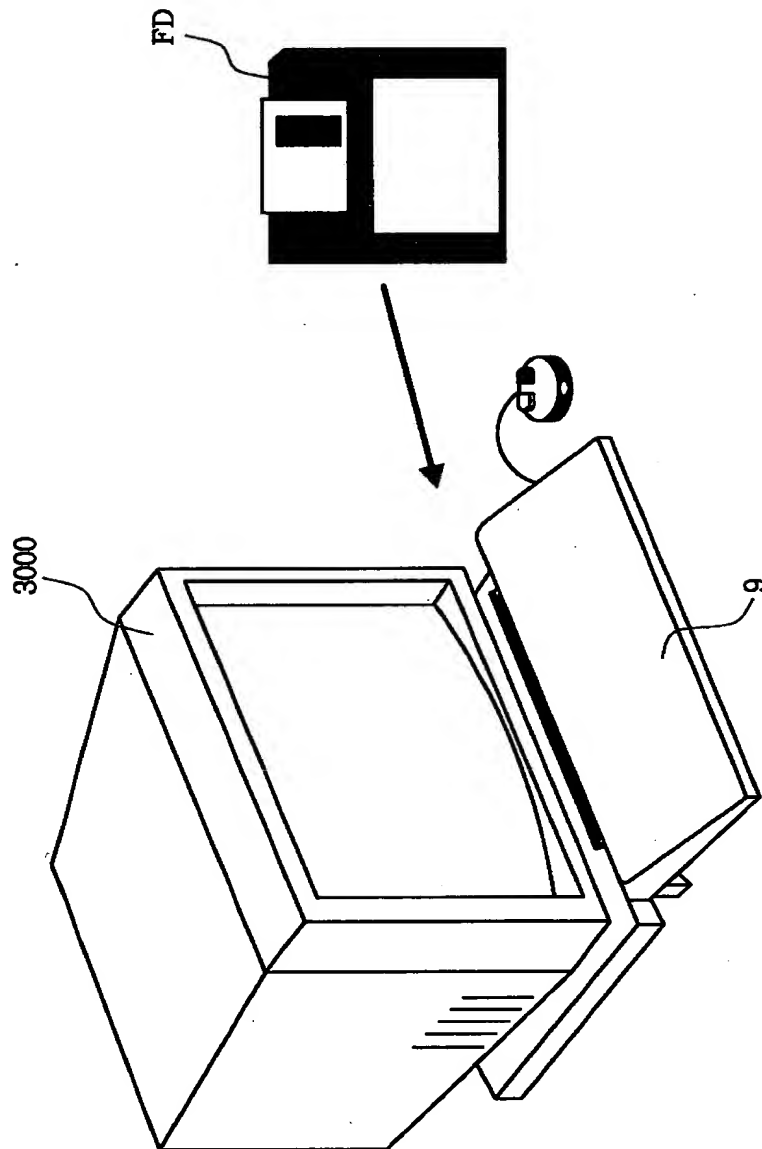
【図 5】



【図6】

1	0.8
2	0.1
3	0.95
4	0.0
5	0.85
6	0.2

【図7】



【書類名】 要約書

【要約】

【課題】 文書等の情報を集合とした集合単位で検索してスコアを求める場合、要素数の多寡によってはスコアが適切とならないことがある。

【解決手段】 文書保持部101の文書から、検索条件保持部103に保持されている検索条件を満足する文書を検索する。検索条件を満足する文書の文書番号を検索結果保持部105に保持する。文書集合毎の文書数と検索結果保持部105中に含まれる文書の数を用いて、集合スコア演算部106により、各文書集合スコアを求めることにより、検索条件により合う文書集合に対して高いスコアが与えられる。

【選択図】 図1

【書類名】 職権訂正データ
【訂正書類】 特許願

<認定情報・付加情報>

【特許出願人】
【識別番号】 000001007
【住所又は居所】 東京都大田区下丸子3丁目30番2号
【氏名又は名称】 キヤノン株式会社
【代理人】 申請人
【識別番号】 100069877
【住所又は居所】 東京都大田区下丸子3-30-2 キヤノン株式会
社内
【氏名又は名称】 丸島 儀一

出 願 人 履 歴 情 報

識別番号 [000001007]

1. 変更年月日	1990年 8月30日
[変更理由]	新規登録
住 所	東京都大田区下丸子3丁目30番2号
氏 名	キヤノン株式会社